# New Statistical Techniques for Predictive Water Quality Modeling at Great Lakes Beaches

Michael N. Fienen[1]

Wesley R. Brooks[1]

Steven R. Corsi[1]

Kurt Wolfe[2]

Rajbir Parmar[2]

Mike Galvin[2]

Mike Cyterski[2]

[1]*USGS Wisconsin Water Science Center, Middleton, WI*
[2]*USEPA Office of Research and Development, Athens, GA*

# Introduction

*E. coli* is used as a fecal indicator bacterium

Without rapid methods, 24 hours are required to get an answer

Environmental predictors (weather, waves, etc.) are available in real time

Our ultimate goal is to predict *E. coli* from environmental predictors.

# Modeling Approaches
*getting the most information from all data available*

**persistence model**

    using yesterday's *E. coli* to predict today's conditions

**ordinary least squares regression (OLS)**

    well-established technique

    requires decisions on how to handle interaction and correlation

**partial least squares (PLS)**

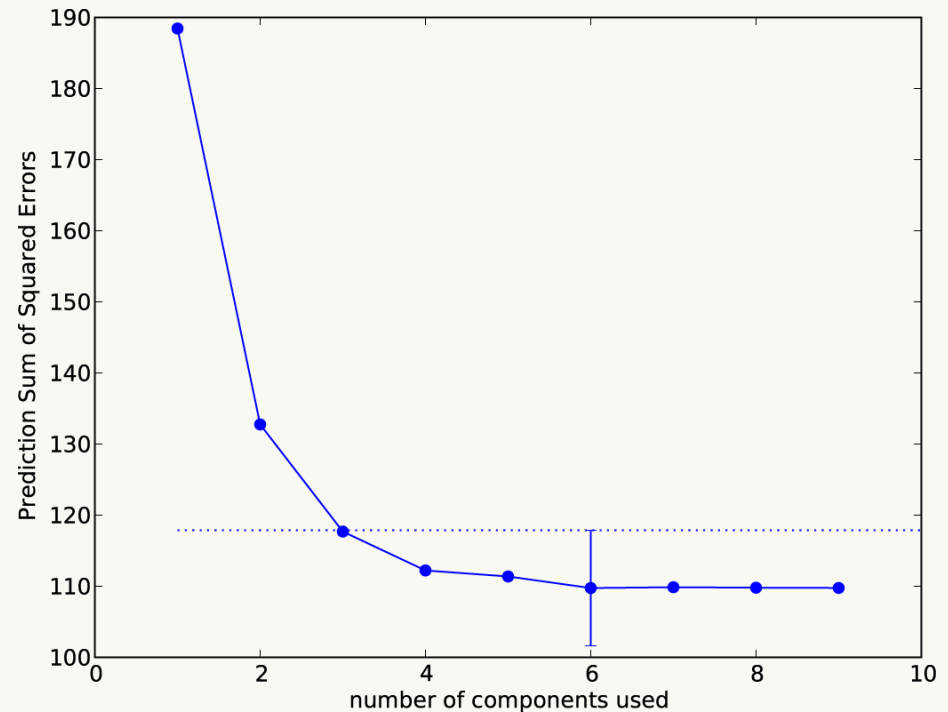    newer technique used extensively in spectroscopy
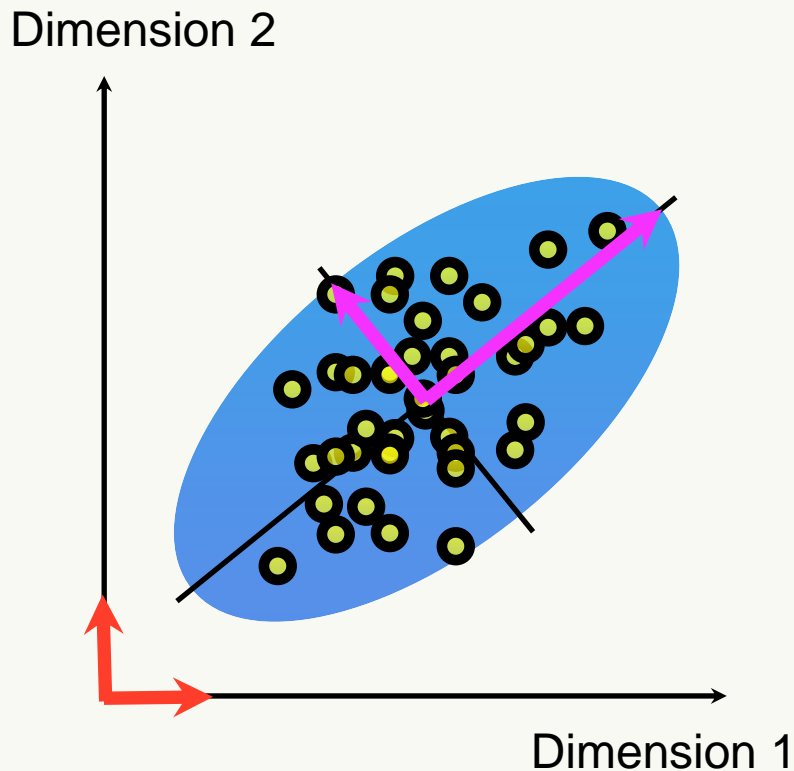
    introduced to beach community by Hou, Rabinovici, and Boehm (2006)

    useful when variables are correlated or insensitive

    overfitting prevented through cross validation and component trimming

    algorithm replaces trial-and-error interaction terms & variable selection

**≋USGS**

# Partial Least Squares: Example

# Partial Least Squares

regression built on principal directions relating variables to predictions

$$\mathbf{y} = \mathbf{X}\beta + \varepsilon \quad\text{—— error term}$$
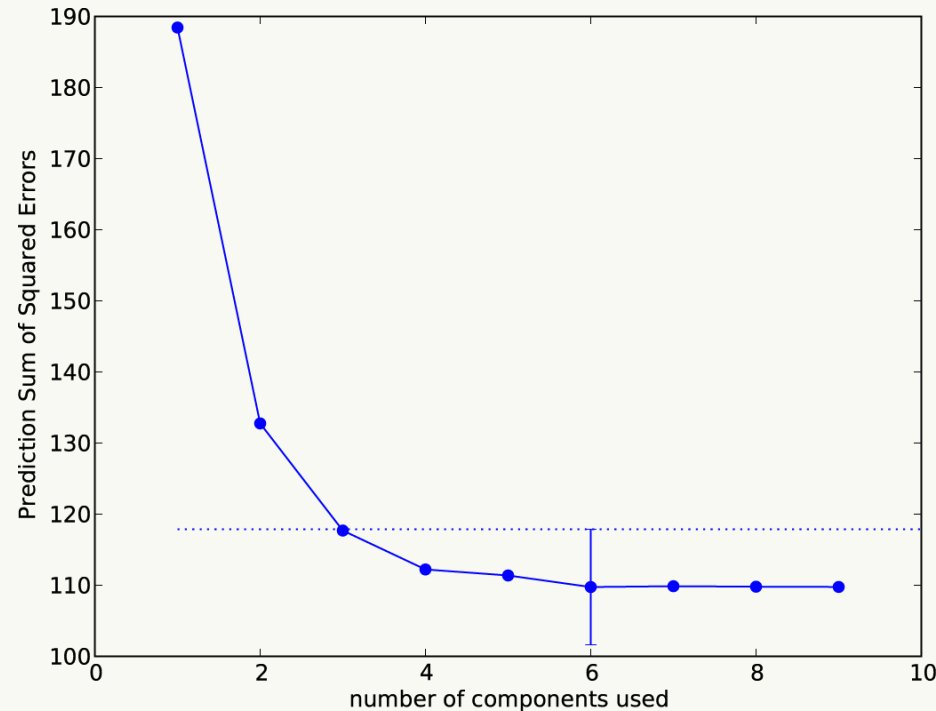
*E. coli* observations

coefficients

predictor variables

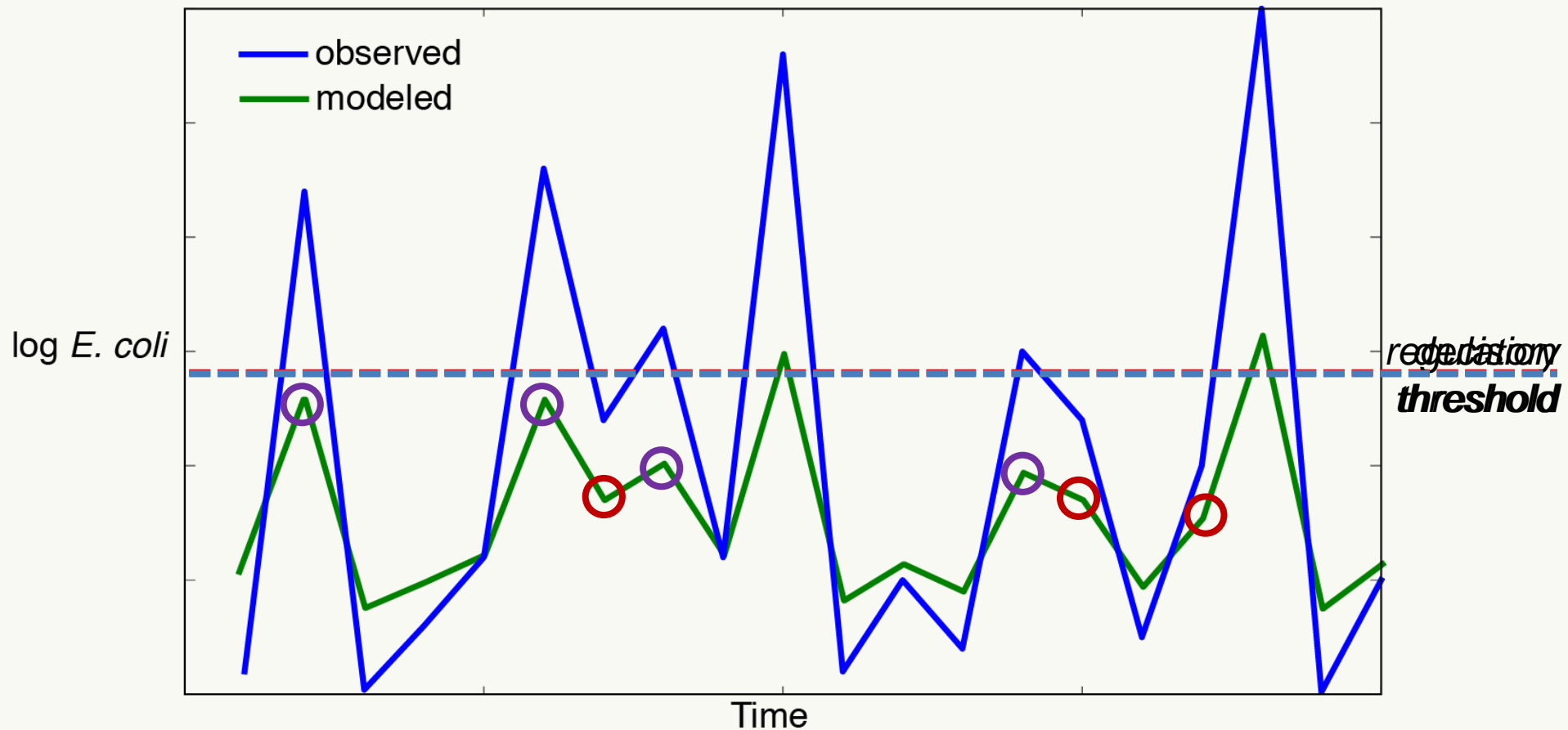principal directions of covariance between $\mathbf{X}$ and $\mathbf{y}$ define components

using all components is equivalent to OLS with all variables

the number of components retained is chosen balancing lower PRESS against overfitting

each component includes information from all base variables
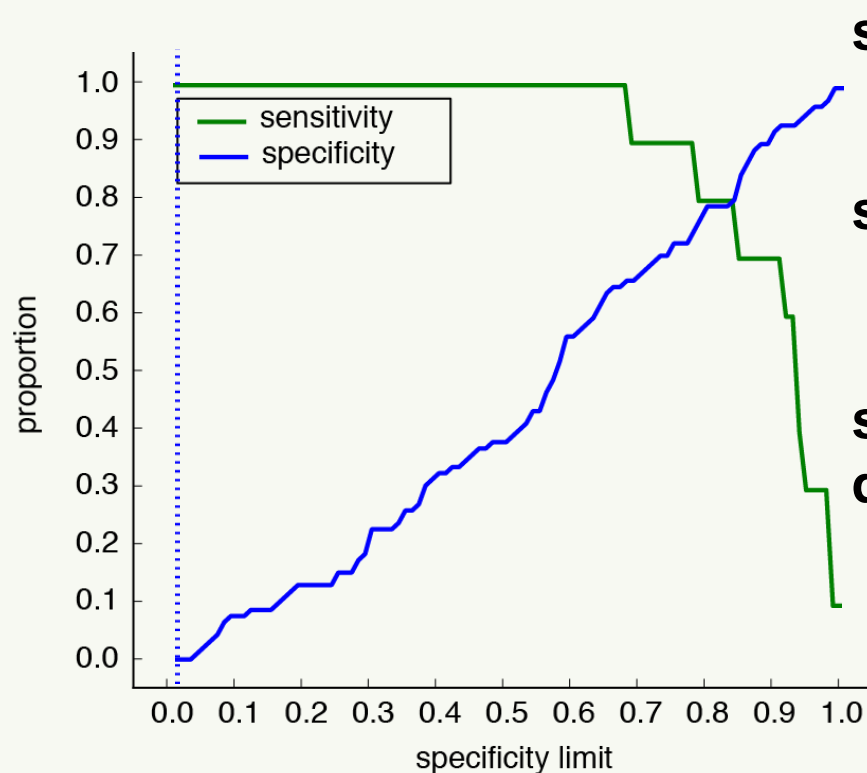


**≋USGS**

# Decision and Regulatory Thresholds

# Tradeoffs for Model Performance
*robust methods and improved data integration*

## building a PLS model

**managers decide the tradeoff between protective and permissive**



**sensitivity: proportion of true positives**
high sensitivity means increasing
true positives

**specificity: proportion of true negatives**
low specificity means increasing
false positives

**specificity limit (**vertical blue line**) is the dial controlling this tradeoff**



≈USGS

# Great Lakes Beaches Modeled



Wisconsin Department of Natural Resources

Upper Lake Park

Maumee Bay State Park

Huntington Reservation

Edgewater State Park

USGS Ohio Water Science Center

USGS

# Building a PLS model – example for Edgewater, Ohio

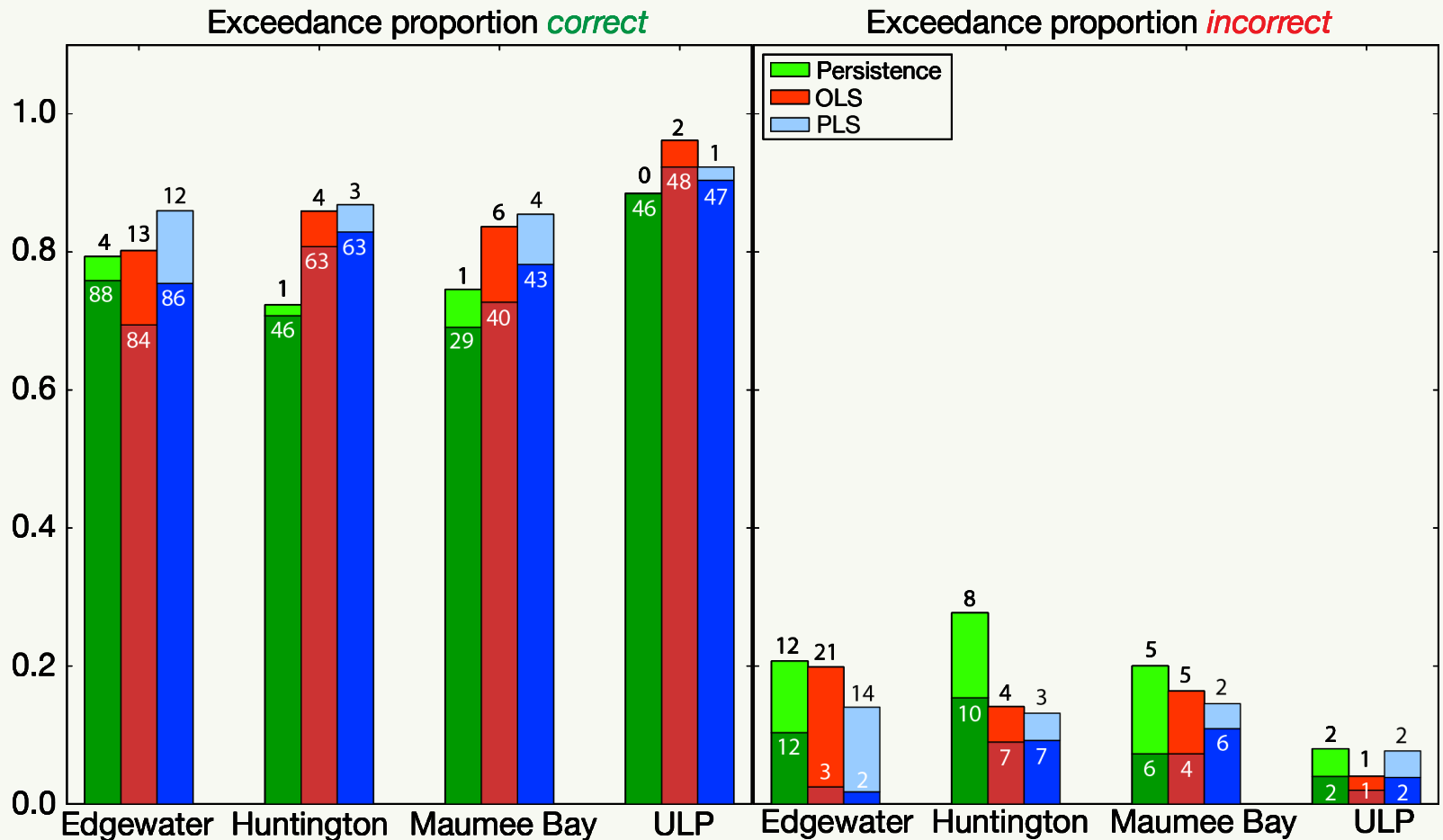| 2005-2010: split training data | Build and test models: predict each fold using data from the other four | 2011: prediction year |
|---|---|---|

**2005-2010: split training data**
- divide the training data randomly into five equally-sized "folds"
- select a few candidate specificity limits (tuning)

**Build and test models:**
- compare model performance on the test folds
- pick the specificity limit that had the best performance, and train a new model over all five folds

**2011: prediction year**
- make and record predictions to manage beach and provide data going forward

≈ **USGS**

# 2010 Model Performance

# Next Steps

**Data acquisition and connectedness**

Greater data acquisition efficiency
and model-building efficiency leads to rapid
application to many beaches

Integration with USEPA and Virtual Beach

# Next Steps

**Virtual Beach**

Well-known EPA software for predicting bacterial concentration

PLS regression to be included in version 3 (coming 2012)

# Acknowledgements



**OLS model results and collaboration:**

Rob Darner

Donna Francy

Amie Brady

Dan Ziegler

Adam Mednick

**Data provided by:**

Cuyahoga County Board of Health

Northeast Ohio Regional Sewer District

Ozaukee County, Wisconsin

**Funding provided by:**

Ocean Research Priorities Plan (ORPP)

USEPA Great Lakes Restoration Initiative (GLRI)

≋ **USGS**